

PATENT

Attorney Docket No. 3373.1

PATENT APPLICATION

METHODS FOR SELECTING NUCLEIC ACID PROBES

Inventors:

Earl Hubbell, A citizen of the United States of America
Residing at 416 S. Genesee
Los Angeles, CA 90036

Assignee:

Affymetrix, Incorporated
A corporation organized under the laws of Delaware

Entity:

Large

Affymetrix, Inc.
3380 Central Expressway
Santa Clara, CA 95051
(408) 731-5000

METHODS FOR SELECTING NUCLEIC ACID PROBES

RELATED APPLICATIONS

5 This application claims the priority of U.S. Provisional Application Number 60/252,617, filed on November 21, 2000, which is incorporated herein by reference for all purposes.

10 This application is related to U.S. Patent Application Number 09/721,042, filed on November 21, 2000, entitled "Methods and Computer Software Products for Predicting Nucleic Acid Hybridization Affinity", and U.S. Patent Application Number 09/718,295, filed on November, 21, 2000, entitled "Methods and Computer Software Products for Selecting Nucleic Acid Probes". Both applications are incorporated herein by reference for all purposes.

BACKGROUND OF THE INVENTION

15 The present invention relates to methods for designing nucleic acid probe arrays. U.S. Patent No. 5,424,186 describes a pioneering technique for, among other things, forming and using high density arrays of molecules such as oligonucleotides, RNA or DNA), peptides, polysaccharides, and other materials. This patent is hereby incorporated
20 by reference for all purposes. However, there is still great need for methods, systems and software for designing high density nucleic acid probe arrays.

SUMMARY OF THE INVENTION

In one aspect of the invention, methods and computer software products are provided for selecting nucleic acid probes. In one embodiment, dynamic programming is employed to select a set of k probes from n probes so that the selected probes have a maximum aggregate adjusted quality score.

A computer implemented method for selecting nucleic acid probes are provided. In some embodiments, the methods include steps of inputting quality scores and locations for a plurality (n) of candidate probes; selecting k number of probes from the n number of candidate probes, wherein the selected probes have a maximum aggregate adjusted quality score; wherein the adjusted quality score is based upon the quality score and the overlapping of the selected probes.

In preferred embodiments, the adjusted quality score is calculated according to:

$$S' = S \sqrt{\frac{l-o}{l}}, \text{ where } S' \text{ is an adjusted quality score; } S \text{ is a quality score; } l \text{ is the}$$

probe length, o is the overlap the probe has with other probes. The probes are particularly useful for measuring gene expression. In preferred embodiments, the probes are immobilized on a substrate to form nucleic acid probe arrays. In exemplary embodiments, the arrays contain probes for measuring a large number of, at least 50, 100, 500, 1000, 2000, 5000 or 10000 transcripts. Each of the transcripts is detected by a set of probes. The embodiment of the invention is often described for the selection of a set of probes for a single transcript (*i.e.*, k probes from n probes), where k is at least 3, 5, 10, or 15.

While this invention is not limited to any particular optimization methods or algorithm, the preferred embodiments employ dynamic programming optimization to select probes. In some particularly preferred embodiments, the selecting step includes calculating best adjusted quality scores ($Score(i, t)$) for probe i last with $t-1$ probes
5 chosen before i and previous location j providing this best score ($Last(i, k)$); determining the best adjusted quality scores for $Score(j, k)$ to select the last probe; and selecting the next probe according to $Last$ (the probe selected, number of probes remain to be selected); and repeating the selecting step until all k probes are selected.

The exemplary embodiments of the systems for selecting nucleic acid probes
10 include a processor; and a memory coupled with the processor, the memory storing a plurality of machine instructions that cause the processor to perform the methods of the invention.

In some embodiments, the computer software products of the invention include a computer readable medium having computer executable instructions for performing the
15 methods of the invention. Exemplary computer readable medium include CD-ROM, DVD-ROM, Floppy Disk, Hard drive, flash memory or the like.

BRIEF DESCRIPTION OF THE DRAWINGS

The accompanying drawings, which are incorporated in and form a part of this
20 specification, illustrate embodiments of the invention and, together with the description, serve to explain the principles of the invention:

Figure 1 illustrates an example of a computer system that may be utilized to execute the software of an embodiment of the invention.

Figure 2 illustrates a system block diagram of the computer system of Figure 1.

Figure 3 illustrates a model system for probe sequence-based prediction of probe hybridization behavior (probe quality).

Figure 4 illustrates a basic physical model for probe target interaction.

Figure 5 shows an example of S_i for an exemplary probe.

Figures 6A and 6B show predicted relative ΔG for perfect match and mismatch probes.

Figure 7 shows an overall reaction of probe target formation.

Figure 8 shows concentration dependency of hybridization intensity.

Figures 9A and 9B show the relationship between probe-target binding affinity (K_{app}) and the slope (S).

Figure 10 shows an embodiment of a process for selecting probes.

Figure 11 shows a pool of candidate probes.

Figure 12 shows another embodiment of a process for selecting probes.

Figure 13 shows yet another embodiment of a process for selecting probes.

Figure 14 shows relationship between overlap of a probe with other probes and the adjustment to quality score in one embodiment of the method of the invention.

Figure 15 shows one embodiment of the computer implemented process of the invention.

Figure 16 shows a process for obtaining weight coefficients.

Figure 17 shows yeast clones used to produce targets.

Figure 18 shows a Latin Square design.

Figure 19 shows Latin Square data sets from yeast_test_hyb chips.

Figure 20 shows a crossvalidation bootstrapping process.

Figures 21A, 21B, 22A and 22B show correlation between predicted and observed hybridization intensities for perfect match probes and mismatch probes.

Figure 23 shows hybridization intensity at different spike concentrations.

Figure 24 shows correlation between predicted and observed intensities over the entire concentration range.

Figure 25 shows predicted versus observed intensities for negative control target.

Figure 26 shows predicted versus observed slopes and improvement of correlation between the two slopes after filtering saturated probes.

Figure 27 shows prediction of hybridization of a human expression chip to human target sequences using weight coefficients generated from the yeast model system.

Figure 28 shows distribution of correlation coefficients.

Figure 29 shows selection of probes using dynamic programming.

Figure 30 compares different methods for selecting probes.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

Reference will now be made in detail to the preferred embodiments of the invention. While the invention will be described in conjunction with the preferred embodiments, it will be understood that they are not intended to limit the invention to these embodiments. On the contrary, the invention is intended to cover alternatives, modifications and equivalents, which may be included within the spirit and scope of the invention.

I. Glossary

“Nucleic acids,” according to the present invention, may include any polymer or oligomer of nucleosides or nucleotides (polynucleotides or oligonucleotides), which include pyrimidine and purine bases, preferably cytosine, thymine, and uracil, and adenine and guanine, respectively. See Albert L. Lehninger, PRINCIPLES OF BIOCHEMISTRY, at 793-800 (Worth Pub. 1982) and L. Stryer BIOCHEMISTRY, 4th Ed., (March 1995), both incorporated by reference. Indeed, the present invention contemplates any deoxyribonucleotide, ribonucleotide or peptide nucleic acid component, and any chemical variants thereof, such as methylated, hydroxymethylated or glucosylated forms of these bases, and the like. The polymers or oligomers may be heterogeneous or homogeneous in composition, and may be isolated from naturally-occurring sources or may be artificially or synthetically produced. See U.S. patent application Serial No. 08/630,427 which is incorporated herein by reference in its entirety for all purposes. In addition, the nucleic acids may be DNA or RNA, or a mixture thereof, and may exist permanently or transitionally in single-stranded or double-stranded

form, including homoduplex, heteroduplex, and hybrid states. Oligonucleotides and polynucleotides are included in this definition and relate to two or more nucleic acids in a polynucleotide.

“Probe,” as used herein, is defined as a nucleic acid, such as an oligonucleotide, capable of binding to a target nucleic acid of complementary sequence through one or more types of chemical bonds, usually through complementary base pairing, usually through hydrogen bond formation. As used herein, a probe may include natural (i.e. A, G, U, C, or T) or modified bases (7-deazaguanosine, inosine, etc.). In addition, the bases in probes may be joined by a linkage other than a phosphodiester bond, so long as the bond does not interfere with hybridization. Thus, probes may be peptide nucleic acids in which the constituent bases are joined by peptide bonds rather than phosphodiester linkages.

“Target nucleic acid” refers to a nucleic acid (often derived from a biological sample), to which the probe is designed to specifically hybridize. It is either the presence or absence of the target nucleic acid that is to be detected, or the amount of the target nucleic acid that is to be quantified. The target nucleic acid has a sequence that is complementary to the nucleic acid sequence of the corresponding probe directed to the target. The term target nucleic acid may refer to the specific subsequence of a larger nucleic acid to which the probe is directed or to the overall sequence (e.g., gene or mRNA) whose expression level it is desired to detect. The difference in usage will be apparent from the context.

An "array" may comprise a solid support with peptide or nucleic acid probes attached to said support. Arrays typically comprise a plurality of different nucleic acids or peptide probes that are coupled to a surface of a substrate in different, known locations. These arrays, also described as "microarrays" or colloquially "chips" have been generally described in the art, for example, U.S. Pat. Nos. 5,143,854, 5,445,934, 5,744,305, 5,677,195, 6,040,193, 5,424,186 and Fodor et al., Science, 251:767-777 (1991). Each of which is incorporated by reference in its entirety for all purposes. These arrays may generally be produced using mechanical synthesis methods or light directed synthesis methods which incorporate a combination of photolithographic methods and solid phase synthesis methods. Techniques for the synthesis of these arrays using mechanical synthesis methods, such as ink jet, channel block, flow channel, and spotting methods which are described in, e.g., U.S. Pat. Nos. 5,384,261, and 6,040,193, which are incorporated herein by reference in their entirety for all purposes. Although a planar array surface is preferred, the array may be fabricated on a surface of virtually any shape or even a multiplicity of surfaces. Arrays may be peptides or nucleic acids on beads, gels, polymeric surfaces, fibers such as fiber optics, glass or any other appropriate substrate, see U.S. Patent Nos. 5,744,305, 5,770,358, 5,789,162, 5,708,153, 6,040,193 and 5,800,992, which are hereby incorporated in their entirety for all purposes. Arrays may be packaged in such a manner as to allow for diagnostics or other manipulation of in an all inclusive device, see for example, US Patent Nos. 5,856,174 and 5,922,591, and 5,945,334, which are incorporated herein in their entirety by reference for all purposes. See also U.S. patent application Serial No. 09/545,207 which is incorporated herein in its

entirety for all purposes for additional information concerning arrays, their manufacture, and their characteristics. It is hereby incorporated by reference in its entirety for all purposes.

II. Probe Selection Systems

5 As will be appreciated by one of skill in the art, the present invention may be embodied as a method, data processing system or program products. Accordingly, the present invention may take the form of data analysis systems, methods, analysis software, etc. Software written according to the present invention is to be stored in some form of computer readable medium, such as memory, or CD-ROM, or transmitted over a
10 network, and executed by a processor. For a description of basic computer systems and computer networks, see, e.g., Introduction to Computing Systems: From Bits and Gates to C and Beyond by Yale N. Patt, Sanjay J. Patel, 1st edition (January 15, 2000) McGraw Hill Text; ISBN: 0072376902; and Introduction to Client/Server Systems : A Practical Guide for Systems Professionals by Paul E. Renaud, 2nd edition (June 1996), John Wiley
15 & Sons; ISBN: 0471133337.

Computer software products may be written in any of various suitable programming languages, such as C, C++, Fortran and Java (Sun Microsystems). The computer software product may be an independent application with data input and data display modules. Alternatively, the computer software products may be classes that may
20 be instantiated as distributed objects. The computer software products may also be component software such as Java Beans (Sun Microsystems), Enterprise Java Beans (EJB), Microsoft® COM/DCOM, etc.

Figure 1 illustrates an example of a computer system that may be used to execute the software of an embodiment of the invention. Figure 1 shows a computer system 101 that includes a display 103, screen 105, cabinet 107, keyboard 109, and mouse 111.

Mouse 111 may have one or more buttons for interacting with a graphic user interface.

5 Cabinet 107 houses a CD-ROM or DVD-ROM drive 102, system memory and a hard drive (113)(*see also* Figure 2) which may be utilized to store and retrieve software programs incorporating computer code that implements the invention, data for use with the invention and the like. Although a CD 114 is shown as an exemplary computer readable medium, other computer readable storage media including floppy disk, tape, 10 flash memory, system memory, and hard drive may be utilized. Additionally, a data signal embodied in a carrier wave (*e.g.*, in a network including the Internet) may be the computer readable storage medium.

Figure 2 shows a system block diagram of computer system 101 used to execute the software of an embodiment of the invention. As in Figure 1, computer system 101 15 includes monitor 103, keyboard 109, and mouse 111. Computer system 101 further includes subsystems such as a central processor 203, system memory 202, fixed storage 210 (*e.g.*, hard drive), removable storage 208 (*e.g.*, CD-ROM), display adapter 206, speakers 204, and network interface 211. Other computer systems suitable for use with the invention may include additional or fewer subsystems. For example, another 20 computer system may include more than one processor 51 or a cache memory. Computer systems suitable for use with the invention may also be embedded in a measurement instrument.

III. Methods for Predicting Quality Scores of Probes

In a preferred embodiment, arrays of oligonucleotides or peptides, for example, are formed on the surface by sequentially removing a photoremovable group from a surface, coupling a monomer to the exposed region of the surface, and repeating the process. These techniques have been used to form extremely dense arrays of oligonucleotides, peptides, and other materials. The synthesis technology associated with this invention has come to be known as "VLSIPS™" or "Very Large Scale Immobilized Polymer Synthesis" technology and is further described below.

Additional techniques for forming and using such arrays are described in U.S. Patent Nos. 5,384,261, and 6,040,193 which are also incorporated by reference for all purposes. Such techniques include systems for mechanically protecting portions of a substrate (or chip), and selectively deprotecting/coupling materials to the substrate. Still further techniques for array synthesis are provided in U.S. Application No. 08/327,512, also incorporated herein by reference for all purposes.

Nucleic acid probe arrays have found wide applications in gene expression monitoring, genotyping and mutation detection. For example, massive parallel gene expression monitoring methods using nucleic acid array technology have been developed to monitor the expression of a large number of genes (e.g., U.S. Patent Numbers 5,871,928, 5,800,992 and 6,040,138; de Saizieu *et al.*, 1998, Bacteria Transcript Imaging by Hybridization of total RNA to Oligonucleotide Arrays, NATURE BIOTECHNOLOGY, 16:45-48; Wodicka *et al.*, 1997, Genome-wide Expression Monitoring in Saccharomyces



cerevisiae, NATURE BIOTECHNOLOGY 15:1359-1367; Lockhart *et al.*, 1996, Expression Monitoring by Hybridization to High Density Oligonucleotide Arrays. NATURE BIOTECHNOLOGY 14:1675-1680; Lander, 1999, Array of Hope, NATURE-GENETICS, 21(suppl.), at 3, all incorporated herein by reference for all purposes). Hybridization-
5 based methodologies for high throughput mutational analysis using high-density oligonucleotide arrays (DNA chips) have been developed, see Hacia *et al.*, 1996, Detection of heterozygous mutations in BRCA1 using high density oligonucleotide arrays and two-color fluorescence analysis. Nat. Genet. 14:441-447, Hacia *et al.*, New approaches to BRCA1 mutation detection, Breast Disease 10:45-59 and Ramsey 1998,
10 DNA chips: State-of-Art, Nat Biotechnol. 16:40-44, all incorporated herein by reference for all purposes). Oligonucleotide arrays have been used to screen for sequence variations in, for example, the *CFTR* gene (U.S. Patent Number 6,027,880, Cronin et al., 1996, Cystic fibrosis mutation detection by hybridization to light-generated DNA probe arrays. Hum. Mut. 7:244-255, both incorporated by reference in their entireties), the
15 human immunodeficiency virus (HIV-1) reverse transcriptase and protease genes (U.S. Patent Number 5,862,242 and Kozal et al., 1996, Extensive polymorphisms observed in HIV-1 clade B protease gene using high density oligonucleotide arrays. Nature Med. 1:735-759, both incorporated herein by reference for all purposes), the mitochondrial genome (Chee et al., 1996, Accessing genetic information with high density DNA arrays. Science 274:610-614) and the BRCA1 gene (U.S. Patent Number 6,013,449, incorporated
20 herein by reference for all purposes).

001237-59654250



In one aspect of the invention, a physical model that is based on the thermodynamic properties of the sequence is used to predict the array-based hybridization intensities of the sequence. Hybridization propensities may be described by energetic parameters derived from the probe sequence, and variations in hybridization and chip manufacturing conditions will result in changes in these parameters that can be detected and corrected. U.S. Patent Application Number 09/721,042, filed November 21, 2000 and incorporated herein by reference, discloses methods for predicting nucleic acid hybridization affinity.

The values of weight coefficients in the physical model may be determined by empirical data because these values are influenced by assay conditions, which include hybridization and target fragmentation, and probe synthesis conditions, which include choice of substrates, coupling efficiency, etc.

In one embodiment (Figure 3), a model experimental system is used to generate empirical data and a computational model is used to process these data to solve for the weight coefficients of the physical model. These solved weight coefficients are in turn placed back into the physical model, enabling it to predict the hybridization behaviors of new sequences.

The interaction between a probe and its target is described in Figure 4. Basically, a target (T) hybridizes to its complementary probe (P) to form a probe-target duplex (P•T) (Figure 4), and the reaction is accompanied with favorable free energy change (Figure 4). The amplitude of the free energy change (ΔG) determines the stability of probe-target duplex. The duplex stability can be described by equilibrium constant

(K_s), which is sequence-dependent. The relationship between K_s and ΔG may be given by Boltzmann's equation:

$$K_s = \frac{k_{on}}{k_{off}} = e^{-\Delta G / RT} \quad [\text{Equation 1}]$$

where k_{on} and k_{off} are the rate constants for association and dissociation, respectively, of

5 the probe-target duplex, R is the gas constant and T is the absolute temperature.

According to Equation 1, ΔG is a function of the sequence. The dependence of ΔG on probe sequence can be quite complicated, but relatively simple models for ΔG have yielded good results.

There are a number of ways to establish the relationship between the sequence and
10 ΔG. In preferred embodiments, one model (equation 2), shown in U.S. Application Serial Number 09/721,042, filed on November 21, 2000, previously incorporated by reference is shown below:

$$\Delta G_{seq} = \sum_{i=1}^{3N} P_i S_i \quad [\text{Equation 2}]$$

or

$$15 \quad \Delta G_{seq} = \sum_{i=1}^{3N} P_i S_i + C \quad [\text{Equation 3}]$$

where N is the length (number of bases) of a probe. P_i is the value of the *i*th parameter which reflects the ΔG of a base in a given sequence position relative to a reference base in the same position. In preferred embodiments, the reference base is A. In this case, the
20 P_i's will be the free energy of a base in a given position relative to base A in the same

position. Figure 5 shows an example of how the value of S_i is determined based upon the sequence of a probe. In this example, a probe, GTCA has $N=4$ and thus, it has $3 \times 4 = 12$ S_i values. Each probe base position has three S values, each for a different possible base. In this example, possible bases are evaluated in the sequence of C, G, and T (A is the reference base). However, one of skill in the art would appreciate that the assignment of this particular base sequence is arbitrary. Alternative evaluation sequence, such as G, C, and T may also be used as long as the same scheme is used for model building and for hybridization affinity prediction.

Based on the simple hybridization scheme described in Figure 4, the hybridization intensity is proportional to the concentration of probe-target duplex, where C_0 is constant (Equation 4). Under equilibrium condition, the intensity is directly related to ΔG (Equation 5). This relationship is also expressed in natural logarithm form, where C_1 and C_2 are constants (Equation 7) and Equation 6 also holds for approaching equilibrium cases. According to Equation 2, the relationship between intensity and probe sequence is described in Equation 7 and 8:

$$I = C_0 [P \cdot T] \quad [\text{Equation 4}]$$

$$[P \cdot T] = K_s [P][T] = e^{-\Delta G/RT} [P][T] \quad [\text{Equation 5}]$$

$$\text{Ln} I = -\Delta G/RT + \text{Ln} \{C_0 [P][T]\} \quad [\text{Equation 6}]$$

C_2

$$\text{Ln} I = C_1 \sum_{i=1}^{3N} P_i S_i + C_2, \text{ where } C_2 =$$

$$\text{Ln} \{C_0 [P][T]\} \text{ and } C_1 = -1/RT \quad [\text{Equation 7}]$$

or

$$\ln I = \sum_{i=1}^{3N} C_i P_i S_i + C_2 = \sum_{i=1}^{3N} W_i S_i + C_2 \quad [\text{Equation 8}]$$

where $W_i = C_i P_i$. The following is a linear regression model for probes of N bases in length using a training data set that contains intensity values of M probes.

$$\ln(I_1) = W_1 S_{11} + W_2 S_{21} + \dots W_{3N} S_{3N1}$$

$$\ln(I_2) = W_1 S_{12} + W_2 S_{22} + \dots W_{3N} S_{3N2}$$

.

.

.

.

$$\ln(I_M) = W_1 S_{1M} + W_2 S_{2M} + \dots W_{3N} S_{3NM}$$

Hybridization intensities (relative to a reference base, such as an A) for each type of bases can be solved at each position in the probe sequence may be predicted. Multiple linear regression analysis is well known in the art, see, for example, the electronic statistic book (<http://www.statsoftinc.com/textbook/stathome.html>); Darlington, R. B. (1990).

Regression and linear models. New York: McGraw-Hill, both incorporated by reference

for all purposes. Computer software packages, such as SAS, SPSS, and MatLib 5.3 provide multiple linear regression functions. In addition, computer software code examples suitable for performing multiple linear regression analysis are provided in, for example, the Numerical Recipes (NR) books developed by Numerical Recipes Software and published by Cambridge University Press (CUP, with U.K. and U.S. web sites).

In a preferred embodiment, a set of probes of different sequences (probes 1 to M) is used as probes in experiments(s). Hybridization affinities (relative ΔG or $\ln(I)$) of the probes with their target are experimentally measured to obtain a training data set (see, example section *infra*). Multiple linear regression may be performed using hybridization



affinities as $I [I_1 \dots I_m]$ to obtain a set of weight coefficients: $[W_1 \dots W_N]$. The weight coefficients are then used to predict the hybridization affinities using Equation 7. Figure 6A shows relative predicted ΔG at every base position in a probe of 25 bases in exemplary experiments (see the example section *infra* for a detailed description of experimental conditions).

In addition, in some embodiments, by using intensities derived from mismatch probes that are probes designed to contain one or more mismatch bases from a reference probe, a set of weight coefficients may be obtained to predict the mismatch intensity using perfect match probe sequence. Figure 6B shows an example for predicting mismatch hybridization affinity at center base position.

Since other interactions such as probe self-folding, probe-to-probe interaction, target self-folding and target-to-target interaction also interfere with the probe-target duplex formation, their contributions to the values of the weight coefficients may also be considered. Figure 7 shows an overall equilibrium scheme including the formation of a probe-target duplex (PT), probe self-folding (P_F) and probe dimerization (PP). Probe folding renders the probe unavailable for binding with the target. Probe dimerization renders two probes unavailable for binding with the target. In some embodiments, the hybridization affinity prediction model accounts for probe folding and probe dimerization:

$$\Delta G_{overall}^0 = -W_d \Delta G_d^0 + W_{PF} \Delta G_{PF}^0 + W_{PP} \Delta G_{pp}^0 \quad [\text{Equation 9}]$$

$$\ln I = C_1 \Delta G_{overall}^0 + C_2 \quad [\text{Equation 10}]$$

where W_d is the weight for sequence based probe affinity; W_{PF} is the weight for probe formation and W_{PP} is the weight for probe dimerization. Any methods that are capable of predicting probe folding and/or probe dimerization are suitable for at least some embodiments of the invention for predicting the hybridization intensity in at least some
5 embodiments of the invention. In a particularly preferred embodiment, Oligowalk (available at <http://rna.chem.rochester.edu/RNAstructure.html>, last visited Nov. 3, 2000) may be used to predict probe folding.

One important criterion of probe selection for a quantitative gene expression assay is that hybridization intensities of the selected probes must correspond to target
10 concentration changes. In some embodiments, the relationship between concentrations and intensities of a probe is modeled as:

$$\ln(I) = S \ln C + \ln K_{app} \quad [\text{Equation 11}]$$

or

$$I = K_{app} C^S \quad [\text{Equation 12}]$$

15 where I is intensity; K_{app} is apparent affinity constant; C is concentration of the target; and S is an empirical value corresponding to the slope of the line relates $\ln I$ and $\ln C$ ($0 < S < 1$) (see Figure 8).

Equation 12 describes the relationship between hybridization intensities of probes and target concentration. For example, when S is equal to 1, the intensities of a probe
20 linearly correspond to its target concentration (Figure 8). Thus, based on the S values of the probe, one can select probes that have good concentration dependence. Figure 9A shows the polynomial relationship between S and $\ln K_{app}$, indicating that when the value

of LnK_{app} increases to a certain level the value of S reaches a plateau before starting to decrease. This relationship allows the identification of not only low hybridization affinity probes (Figure 9B, bottom lines) but also GC-rich probes that have high affinity but bind to both specific and non-specific targets (Figure 9B, top line). These GC-rich probes have high intensities, but the intensities maintain constant when target concentration changes (Figure 9B, top line). Therefore, these probes have small slopes. In some embodiments, linear regression modeling alone will not identify probes with a high propensity to saturate. That is because the linear model for each target concentration will predict the intensity that a probe would have had if it could bind to an unlimited amount of target. Therefore, the predicted slope can be quite high when the observed slope is low (Figure 26, top). The well-behaved relationship between predicted LnK_{app} and observed slope allows filtering probes with a high propensity to saturate based on the predicted LnK_{app} for the given probe. If LnK_{app} is above a cutoff value (e.g., 5, 6, 7, 8, 9, or 10, Figure 9), then the probe is effectively filtered as a candidate for probe selection. Figure 26 (middle) shows the predicted slope profiles after filtering as well as the significant improvement in the overall correlation after these regions are removed.

IV. Methods and Software for Selecting Probes

Figure 10 shows a computer-implemented process for selecting probe sequences from a pool of candidate probes. In this particular embodiment, the sequences of a pool of candidate oligonucleotide probe are processed by a quality predictor (101). Throughout this application, the term probe may refer to the sequence of a probe. The

09745965-122100

pool of candidate oligonucleotide probes may be all possible probes against a particular target or targets. Typically, oligonucleotide probes are at least 10, 15, 20, 25 and 30 bases in length. Polynucleotide probes can be more than 10, 20, 25, 30, 100, 200, 500, 1000, or 5000 bases in length. Figure 11 illustrates a complete pool of candidate oligonucleotide probes (unfiled rectangle boxes) against a target (black rectangle box). Each of the probes is designed to be complementary to the target sequence. In this particular embodiment, the oligonucleotides are 25mers. The first probe is complementary to bases 1-25 (from the 5' end) of the target sequence. The second probe is complementary to bases 2-26 and so on. While a complete pool is often desirable, it is not necessary to have a complete pool for at least some embodiments of the invention. In some cases, filters may be used to remove some of the probes from the pool.

The input to the quality predictor (Figure 10, 101) is the sequence of a pool of candidate probes. One of skill in the art would appreciate that the format of input is not critical. In some embodiments, the probe sequences may be inputted from one or a number of probe sequence files. The file(s) may be plain text file(s), in the FASTA format or other suitable file format. Alternatively, the input may be a stream from other sources such as a data pocket stream from a remote networked computer.

The quality predicator is a software module that calculates quality scores (the term score refers to any qualitative and quantitative values with regard to desired properties of a probe) for probes based upon the sequences of probes. In some embodiments, the quality score may include predicted values such as perfect match intensity, mismatch intensity and/or slope.

Probe selection module (103) selects probes based upon their scores. In preferred embodiments, the quality scores are combined to obtain a unified score. In some cases, the unified quality score is the simple summation of quality scores (e.g., Unified Quality Score=Perfect Match Intensity+Mismatch Intensity+Slope). The selection of probes may be based upon the scores only. For example, if a certain number of probes are desired, the probes with the highest scores are selected until enough number of probes are selected. Alternatively, a threshold-unified score may be established. Probes that have scores higher than the threshold score are selected.

In a preferred embodiment, the goal of the probe selection step is to find the best probes to represent a sequence. The probe selection software module takes a set of probes and a set of quality measures for each probe. It then implements an optimization algorithm to find the best n probes, spread out across the gene. Methods for probe selection using optimization algorithm is described in U.S. Application Number 60/252,617, filed November 21, 2000, and incorporated herein by reference in its entirety for all purposes.

Figure 12 shows another embodiment of the computer implemented probe selection process of the invention. Target sequences are inputted to a candidate probe generator (121) which produce either all possible probes of certain length or a subset of the all possible probes. The candidate probe sequences are fed to the quality score predictor (122) for calculating quality measures (scores, e.g., perfect match intensity, mismatch intensity and/or slope). The candidate probe sequences are also fed to a 3' bias score predictor (123) to obtain 3' bias scores that indicates the distance of probe sequence

from the 3' end of target sequence. Since the current target preparation method is 3' biased, it is important to select probes that fall into range where its target will be made. The probe sequences may optionally be inputted into a cross hybridization score predictor (124) to calculate cross hybridization scores. The quality scores, 3' bias scores and/or cross hybridization score are combined by a probe score calculator module (125) to produce a unified score. A probe selection module (126) picks the probes with the desirable score.

Figure 13 shows a complete computer implemented probe selection process. In this preferred embodiment, target sequences (131) are used to generate a pool of candidate probes. The probe sequences are stored in a FASTA sequence file. A sequence file splitter (132) divides probe sequences to .seq file which store one sequence per file. The .seq files are processed using an OligoWalker batch tool (133) to produce a .rep file, one for each probe sequence. The .rep files contain ΔG values for the probes. The rep files are inputted into a quality predictor (134). The quality predictor is based upon multiple linear regression models derived from experiment data using, for example, yeast test chips (see also, example section below) (1310). The quality predictor calculates quality scores (measures, perfect match intensity, mismatch intensity and slope) as described above in section II. The rep file is also inputted into a 3' bias score predictor (135) to estimate 3' bias scores for the probes.

The multiple probe FASTA sequence file is also inputted into a cross hybridization predictor (136) to predict a cross hybridization score. The cross hybridization score predictor is based upon models (such as multiple linear regression

models) derived from experiment data (1311). In some embodiments, cross hybridization may also be evaluated by pruning probe sequences against a human genome data base (1312) which may be residing locally, in a local area network or in a remote site such as the Genbank (<http://www.ncbi.nlm.nih.gov>).

- 5 The quality measures, 3' bias scores and cross hybridization scores are combined by the probe score calculator (137) to produce a unified score for each probe. The combined score is then used for selecting probes (138). The probe selection module takes a set of probes and a set of quality measures for each probe. It then implements a dynamic programming algorithm to find the best n probes, spread out across the gene.
- 10 The selected probe sequences are stored in .101 files (139).

The following tables describe the various software modules in the exemplary embodiments described in Figure 13.

1: Multiple linear regression modeling tool

Description	Calculates the weights for the regression model. Its is a one time calculation. The results of the calculations will be used every time a new chip is designed.
Input	Yeast Test Chip, available from Affymetrix, Santa Clara, CA
Output	Multiple linear regression models, a set of weights.
Part of chip design	In this embodiment, it is not part of the software package for chip design. It is used as a one-time external process. However, in other exemplary embodiments, it may also become part of the software.

2: Sequence file splitter

Description	Splits a FASTA file of sequences into several sequence files one for each sequence in the instruction file. If max files in folder is greater than 0, subfolders are created in the output path. Each subfolder gets up to the maximum files specified.
Input	<ul style="list-style-type: none">◦ FASTA file◦ Instructions file◦ Output path◦ Max files in one folder
Language / Tool	Java

3: Oligo Walk batch tool

Description	Runs Oligo Walk in batch mode. Oligo Walk produces a .rep file for each sequence. The .rep file contains a delta G value for each probe
Input	Batch of .seq files
Output	.rep file. The .rep file identifies a probe by a number and a sequence. The sequence is a reverse complement of the 25-mer it represents on the input sequence. The number is the beginning of the probe.
Part of chip design	Yes

Language / Tool	Microsoft ® Visual Basic
------------------------	--------------------------

4: Quality predictor

Description	Takes in the MLR model measures and delta G values from Oligo Walk and produces 3 quality measures, perfect match intensity, mismatch intensity and slope.
Input	.rep file produced from Oligo Walk
Output	3 Quality measures for each probe. The probe is described as in the input format.
Part of chip design	Yes
Language / Tool	C

5: Cross Hyb Modeling Tool

Description	Analyzes the results of the yeast cross hyb chip to create a model for predicting the cross hyb score for a probe, based on the number of mismatches and positions of mismatches with 1 or more matching sequences.
Input	Results from the cross hyb chip
Output	A model that relates number of mismatches and positions of mismatches to a cross hyb score.
Part of chip	In some embodiments, it is not part of the chip design package.



design	Alternatively, it can be part of the package.
---------------	---

6: Cross Hyb Score Predictor

Description	Predicts a cross hyb score for a given set of probes. Its does so by matching the given probes with a genome and assigns a numeric score using the cross hyb models.
Input	<ul style="list-style-type: none">▫ Cross hyb models▫ A genome▫ Set of probes
Output	List of probes and corresponding cross hyb scores
Part of chip design	Yes

5 7: 3' Bias Score Predictor

Description	Predicts the 3' bias score for a given set of probes. Earlier it was believed that most sequences have a sigmoid graph for the 3' bias. But, recently used sequences do not always follow the pattern Therefore, it is important to first study the 3' bias effect and then design a measurement model.
Input	<ul style="list-style-type: none">▫ Set of probes
Output	List of probes and corresponding 3' bias scores

Part of chip design	Yes
----------------------------	-----

8: Probe Score Calculator

Description	Given a set files with probe information and scores, this program matches each probe in each sequence and calculates 1 unified score for each probe.
Input	<ul style="list-style-type: none"> ▫ Set of probes ▫ Set of measures for each probe, each in a different file(s) <ul style="list-style-type: none"> ▫ 3 quality scores (probes defined in OligoWalk format) ▫ cross hyb score (probes defined in chip design format) ▫ 3' bias score (probes defined in chip design format)
Output	List of probes with a corresponding score.
Part of chip design	Yes

9: Probe selection algorithm

Description	Finds the best probes to represent a sequence. It takes a set of probes and a set of quality measures for each probe. It then implements a dynamic programming algorithm to find the best n probes, spread out across the gene.
Input	<ul style="list-style-type: none"> ▫ Set of probes ▫ Set of measures for each probe <ul style="list-style-type: none"> ▫ 3 quality scores ▫ cross hyb score ▫ 3' bias score ▫ Number of probes to choose
Output	.llq file
Part of chip design	Yes
Language / Tool	C

10: Algorithm Test Tool

Description	Tests the new probe selection algorithm. The probe selection algorithm is used to select probes for the known, Yeast test chip. The selected probes are analyzed for their intensity, slope and discrimination values on the yeast test chip.
Input	<ul style="list-style-type: none"> ▫ Probes selected for the sequences on the yeast test chip



	▪ Results from the yeast test chip
--	------------------------------------

V. Dynamic Programming for Probe Selection

When DNA arrays are constructed, it is vitally important to choose the best set of probes for the type of analysis that will be done. In particular, for any particular application, it is possible to assign scores to the probes (such as the quality score described above), so that probes with higher scores are more likely to be better suited for a particular application than others. Given a set of probes with scores, it is desirable to pick the best set of probes.

In selecting multiple nucleic acid probes for one target, one complication that arises is that probes that are nearby each other are mostly redundant. The amount of new data observed from a probe that overlaps with another probe by 24 bases out of 25 is minimal, so that even if such a probe has a high quality score, it may be desirable to pick another probe that has a lower quality score, but has no or less overlaps with other probes.

In one aspect of the invention, methods are provided to adjust the quality score for each probe corresponding to the amount of information it would provide. In some embodiments of the methods, the following rules are used to adjust the score:

- 1) Probes that do not overlap have full scores;
- 2) Probes that do overlap have a penalty that decreases as the amount of penalty decreases;
- 3) Scores are correlated with information provided;
- 4) Adding scores provides a reasonable estimate of “total information”

5) It is only necessary to consider overlap with the previous probe for estimating new information.

In particularly preferred embodiments, the quality score is adjusted as follows:

$$S' = S \sqrt{\frac{l-o}{l}}, \text{ where } S' \text{ is adjusted score; } S \text{ is initial score; } l \text{ is the probe length, } o \text{ is the}$$

5 overlap the probe has with other probes. Figure 14 shows the relationship between overlapping bases and the penalty that may be applied when the probes are of 25 bases in length. In the preferred embodiments, if a probe has no overlap over other probe(s), the adjustment is 1.0 or no penalty. The penalty increases as the number of overlapping bases increases. For example, if a probe has 10 overlapping bases, the adjustment quality
10 score is about 77% of the original quality score. For a probe with 24 bases overlapping with other probes, the original quality score is adjusted for a 80% penalty because of the additional information content provided by this probe in view of the other probes is small.

One of skill in the art would appreciate that the methods of the invention are not limited to any particular methods for adjusting quality scores.

15 In one aspect of the invention, optimization methods are used to pick an optimal set of k probes from n probes provided with initial scores and locations of the probes in the target sequence. The optimal set of k probes is chosen for its high (optimal) aggregate, not individual adjusted score. In a typical gene expression experiment, k may be at least 3, 5, 10, 15, 20, 25, 30, 40 or 50 for a single transcript. The selection process
20 may be described with reference to a single transcript and the selection of a single set. The methods are particularly useful for selecting probes against a large number of

transcripts, for example, at least 100, 200, 300, 500, 1000, 2000, 5000, or 10000. A set of probes may be selected for each of the transcripts.

In some embodiments, dynamic programming is used to select the optimal set of probes with maximum aggregate adjusted scores. Figure 15 shows an exemplary

5 computer implemented process.

A computer program starts (1501) and inputs (1502) quality scores ($score(i)$) and location of the probes in the target sequence.

Step 1503 calculates $Score(i,t)$, i.e., best score using probe i last with $t-1$ probes chosen before i and $Last(i,t)$, i.e., previous location j providing this best score. The

10 following is an exemplary pseudo-code for this process:

CalculatingBestScore() {

//Initializing - the score for the first probe selected

for ($i=1$; $i \leq n$; $i++$) {

Score($i,1$) = $score(i)$;

Last($i,1$) = I ;

}

//running through all scores for selecting additional probes

for ($t=2$; $t \leq k$; $t++$) {

for ($i=1$; $i \leq n$; $i++$) {

Score(i,t)

$$= \max \{ j \text{ in } 1 \dots i-1 \mid \text{Score}(j, t-1) = \text{score}(i) * \text{Screen}(\text{position}(i) - \text{position}(j)) \}$$

$$\text{Last}(i, t) = j; \text{ //last } (i, j) \text{ records where the max occurs}$$

}

5 }

}

This algorithm may be accelerated by utilizing the fact that probes that do not overlap have full scores, so not all 'j' have to be searched over.

The best set of k probes given the scores can be found by backtracking through the Score matrix to extract the k probes that together yield the best score (Steps 1504, 1505, 1506, 1507). Step 1504 finds the best score for $\text{Score}(i, k)$. The "last" probe selected is probe i . $j = \text{Last}(i, k)$ at this location gives us the next-to-last probe selected for the best set (1504). Similarly, $\text{Last}(j, k-1)$ gives us the next-to-next-to-last probe (1505).

The following is an exemplary pseudo-code for this process:

```

15       FindBestProbes(){
           l[0] = bloc; // last score added, l[0] is the location of the last probe selected
           i = Last(l[0], k); //i the last probe selected.
           for (g=k-1; g>=1; g--) {
               p = Last(i,g);
20               l[k-g] = i;
               i = k;
           }

```

}
In some embodiments, an additive gap score penalty (as opposed to the multiplicative described above) is used, but it seems that the multiplicative penalty provides better results. Particularly preferred embodiments employ dynamic programming to select probes. The dynamic program optimizes the best set of probes, rather than optimizing individual probes. The gap penalty formulation is very flexible, and allows for explicit trade-offs between distance and quality. Because the gap penalty stops changing after a certain distance, the algorithm may be accelerated, and run in time proportional to $k*n*(length\ of\ penalty)$, much, much faster than $k*n*n$ without acceleration.

VI. Examples

The following examples demonstrate the effectiveness of the methods of the invention for predicting hybridization intensities and for selecting oligonucleotide probes for gene expression monitoring.

A. Example 1: Prediction of Hybridization Intensities of Probes Against Yeast Genes

Figure 16 shows the overall process of the experiments. Yeast was used as a model system for this experiment because the yeast genome had been sequenced. Arrays containing nucleic acid probes complementary to yeast genes are commercially available from Affymetrix (Santa Clara, California). Genes were selected to cover sequence complexity such as GC content, secondary structure, Motif and gene clusters. Twenty

probe pairs (perfect match and mismatch probes) were selected to cover the entire sequence of one of the 112 selected yeast genes. The probes are synthesized in situ on glass substrate using photo-directed synthesis method that was disclosed in, for example, U.S. Patent Nos. 5,384,261, 5,744,305, 5,445,934 and 6,040,138.

One hundred and twelve yeast clones representing the 112 genes were randomly divided into 14 groups (Figure 17). Labeled targets prepared from these clones were used as spikes for 14 experiments at various concentration levels from 0pM to 1024pM. In some experiments, the spikes derived from yeast gene clones were combined with labeled nucleic acid representing human complex background. A 14 x 14 Latin square design (Figure 18) was employed. The numbers in the table indicates the concentration used (pM). For each experiment, 14 groups of genes at 14 different concentrations were pooled together and hybridized to an oligonucleotide probe array. For each Latin Square 14 oligonucleotide probe array hybridization experiments were performed. Figure 19 shows experiments conducted.

Cross-validation (Figure 20) was used to evaluate the prediction. The cross-validation process held one gene for test and used the other 111 genes to solve the weight coefficients that in turn were used to predict intensities for the test genes, as described in Figure 20. The correlation between the predicted and measured intensity for one test gene (YDR113C) is shown in Figure 21A, and Figure 21B shows the correlation against target sequence, where lines represent the predicted values and dots represent the observed values. The correlation of the predicted and measured values for perfect match (PM) and mismatch (MM) probes is also demonstrated in Figures 22A and 22B respectively, where

lines represent the predicted values and dots represent the observed values for gene YGR109C.

Figure 23 shows predicted intensity versus actual intensity at various target spike concentrations, where lines indicate the predicted values and dots represent observed values. Figure 24 shows correlation coefficients between predicted and observed intensity (LnI) as function of concentration, where top and bottom lines represent perfect matches and mismatches, respectively. The high correlation (0.85) holds for 4000-fold concentration range (Figure 24), and the results demonstrate that the methods of invention are able to predict probe behaviors through a wide dynamic range.

Figure 25 shows predicted versus observed intensities when the target transcripts were derived from genes in the wrong orientation, which resulted in no complimentary target generated for the probes. As shown in Figure 25, predicted intensities (lines) had no correlation with observed intensity (dots) because right target is absent. The result indicates the prediction method is accurate and specific.

Figure 26 shows predicted slope versus observed slope. In some regions in Figure 26 (top), the values of predicted slope (lines) can be quite high when the values of observed slope (dots) because of the saturated probes in those regions. According to Equation 12 and Figure 9, the saturated probes can be identified and removed. Figure 26 (middle) shows the predicted slope profiles after filtering the saturated probes and the significant improvement in the overall correlation after these regions are removed.

B. Example 2: Prediction of Hybridization Intensities of Probes from Human Genes

This example demonstrates that weight coefficients obtained from the model yeast experiment system is also able to predict the intensities on the human gene expression chip and the predicted intensities (left bar) are highly correlated with observed intensities (right bar) at each probe position as indicated by x-axis. The correlation is shown in

5 Figures 27 A-E. Typically, the correlation coefficients ranged from 0.45-0.83. The distribution of the correlation coefficients are shown in Figure 28. These results demonstrate that the probe selection model may be generalized to different organisms such as mammals, plants.

10 C. Example 3: Probe Selection

This example demonstrates that the model-based probe selection method and software may provide improvement over current probe selection methods. Figure 29 shows intensity values of sixteen probes (open squares) selected for the Yer161c gene based upon quality scores and using dynamic programming. Figure 29 also shows that

15 the sixteen selected probes (open squares) are spaced along sequence. Figure 30 shows a comparison of average intensity difference (between perfect match and mismatch) values of probe selected by various methods for all yeast test genes. Probes selected randomly (diamonds) were similar to those selected according empirical rules (squares). The model based selection method (triangles) improved average intensity difference values. The

20 result indicates the model-selected probes have high sensitivity and specificity.

Conclusion

The present invention provides methods and computer software products for predicting nucleic acid hybridization affinity, detecting mutation, selecting better-behaved probes, and improving probe array manufacturing quality control. It is to be understood that the above description is intended to be illustrative and not restrictive. Many variations of the invention will be apparent to those of skill in the art upon reviewing the above description. By way of example, the invention has been described primarily with reference to the design of a high density oligonucleotide array, but it will be readily recognized by those of skill in the art that the methods may be used to predict the hybridization affinity of other immobilized probes, such as probes that are immobilized in or on optical fibers or other supports by any deposition methods. The basic methods and computer software of the invention may also be used to predict solution-based hybridization. The scope of the invention should, therefore, be determined with reference to the appended claims, along with the full scope of equivalents to which such claims are entitled. All references cited herein are incorporated herein by reference for all purposes.